

# Credit Default Prediction

## Machine Learning Final Project

Group 5 : Sameer Batra & Jeongmin An

# Table of Contents

## Overview

- Project Goals
- Dataset overview

## Data Pre-processing

- Handling Missing Value
- Handling Outliers

## Exploratory Data Analysis

- Target Variable Analysis
- Age distribution
- Monthly Income Analysis
- Late payment Behavior
- Feature Correlation Analysis

## Hyperparameter Tuning & Model selection

# Project Goals



## **Risk Assessment**

- Predict future delinquency to strengthen risk monitoring
- Reduce financial losses by identifying high-risk borrowers early



## **Support Data-Driven Credit Decisions**

- Provide consistent and responsible lending decisions
- Deliver actionable insights into financial behaviors and risk drivers



## **Build a reliable machine learning model**

- Compare multiple models and select the best-performing approach
- Ensure strong performance through preprocessing, tuning, and evaluation

# Dataset Overview

## “Give Me Some Credit” Dataset

- ~150,000 records from Kaggle Competition
- The target variable, **SeriousDlqin2yrs**, indicates whether a person becomes 90+ days delinquent **within the two years**, making this a binary classification problem.

The dataset contains **11 variables**, including key financial indicators such as:

**RevolvingUtilizationOfUnsecuredLines** (credit usage rate)

**DebtRatio**

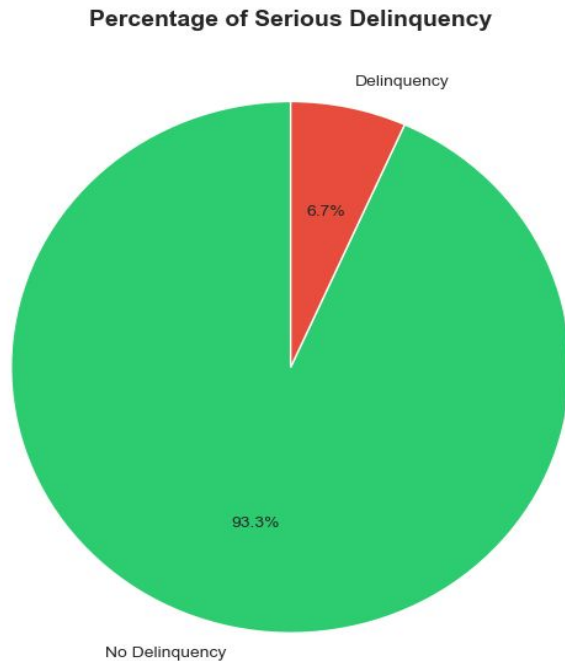
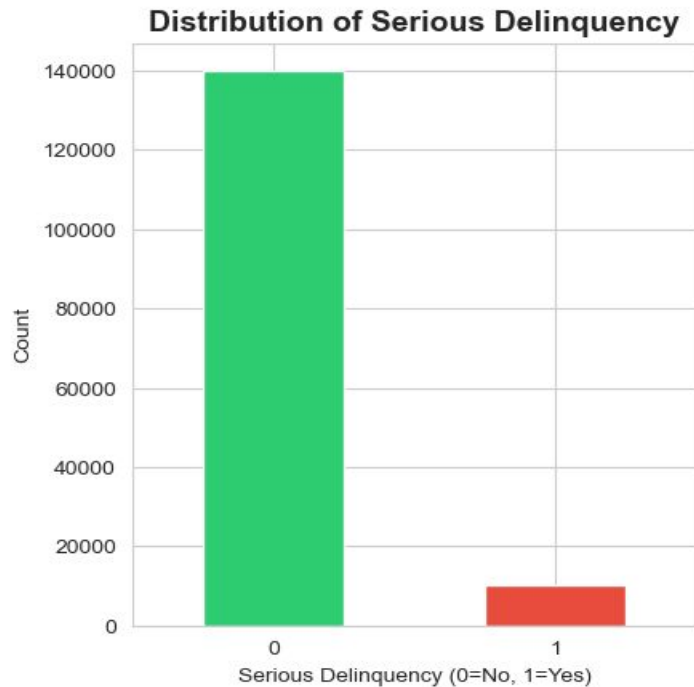
**MonthlyIncome**

**Delinquency history**  
at 30–59, 60–89,  
and 90+ days

- **Number of open credit lines**
- **Real estate loans**
- **Dependents**

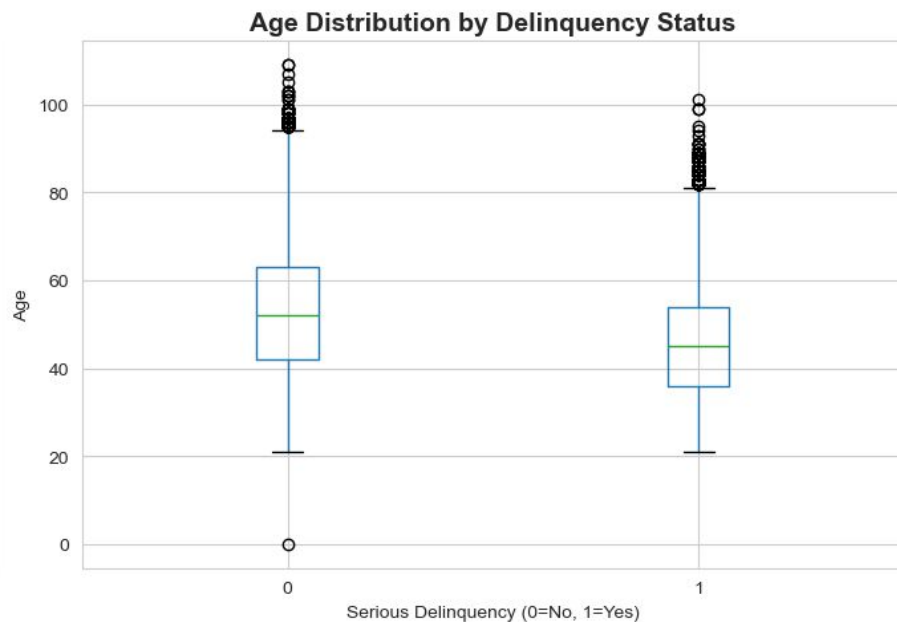
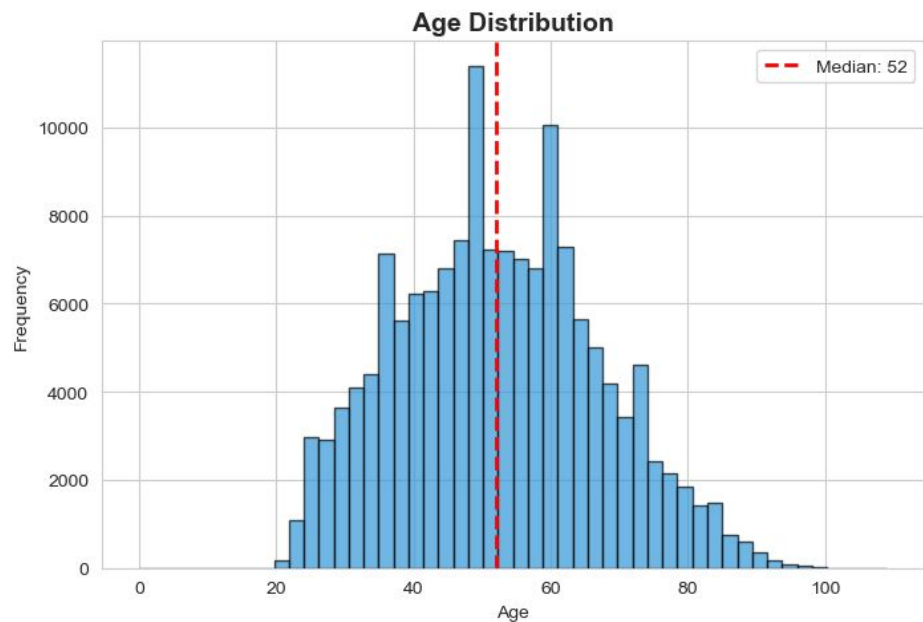
# EDA

## Target Variable Analysis



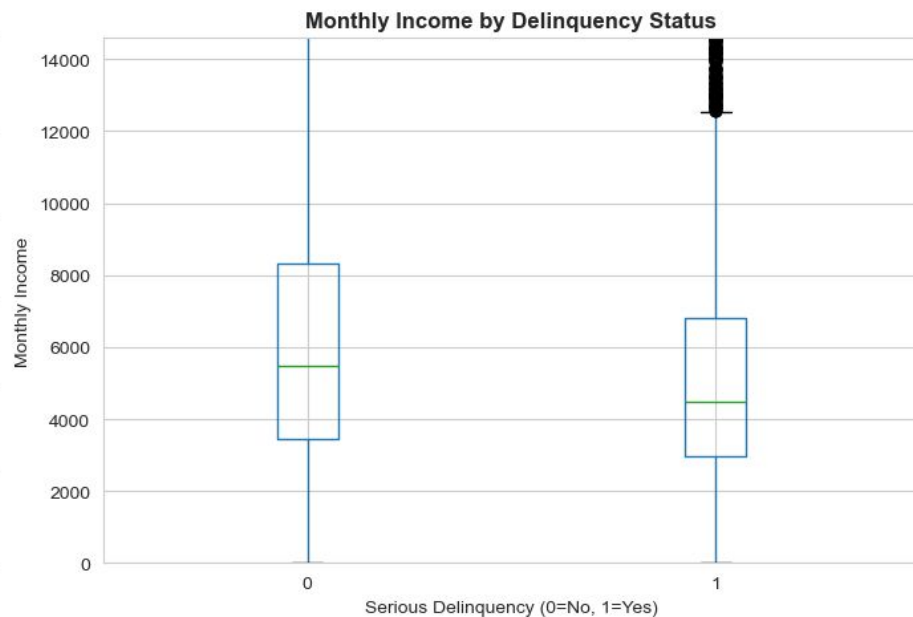
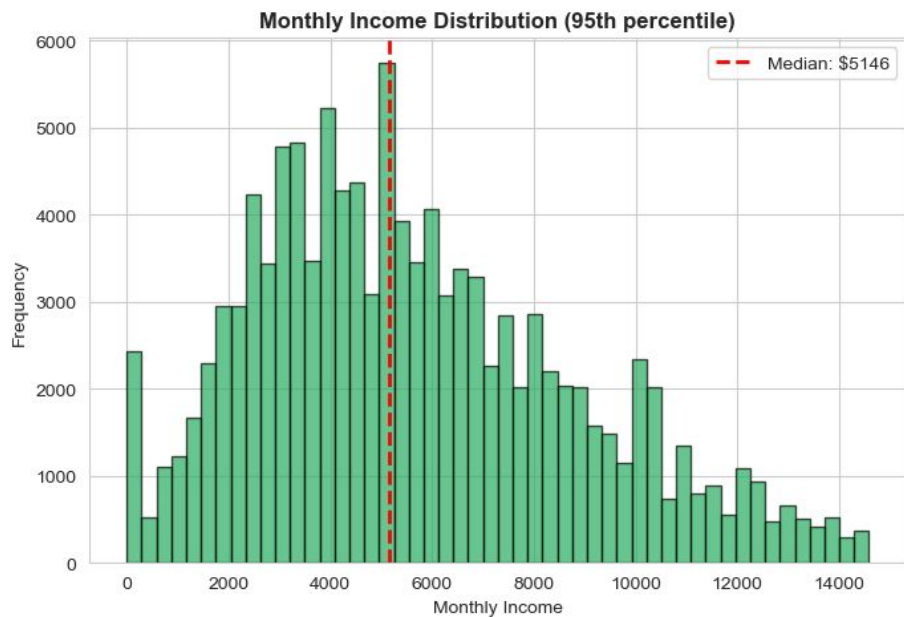
# EDA

## Age Distribution



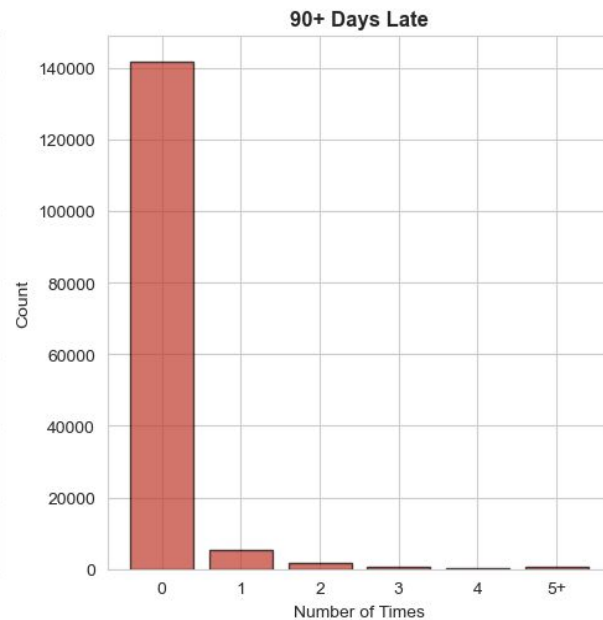
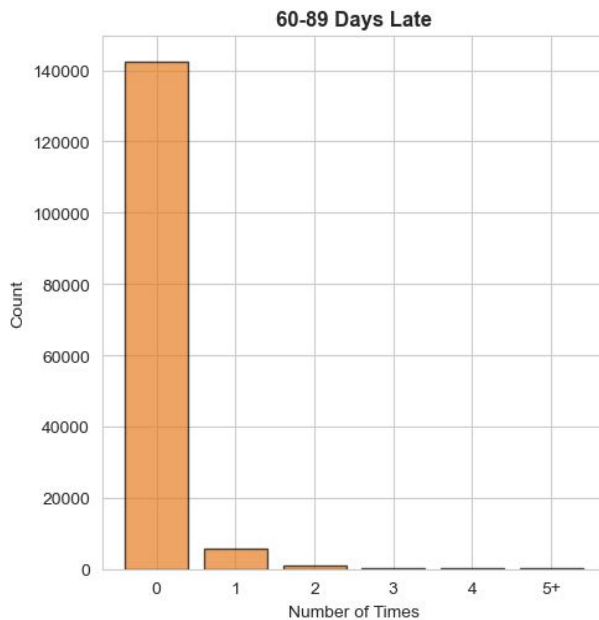
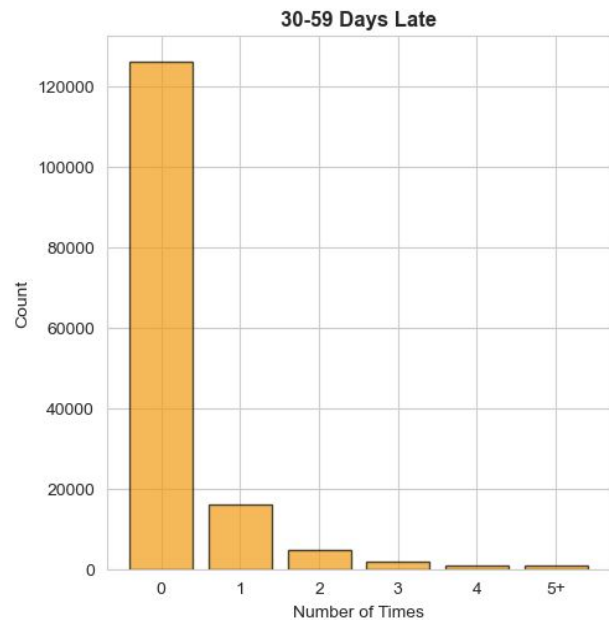
# EDA

## Monthly Income Analysis



# EDA

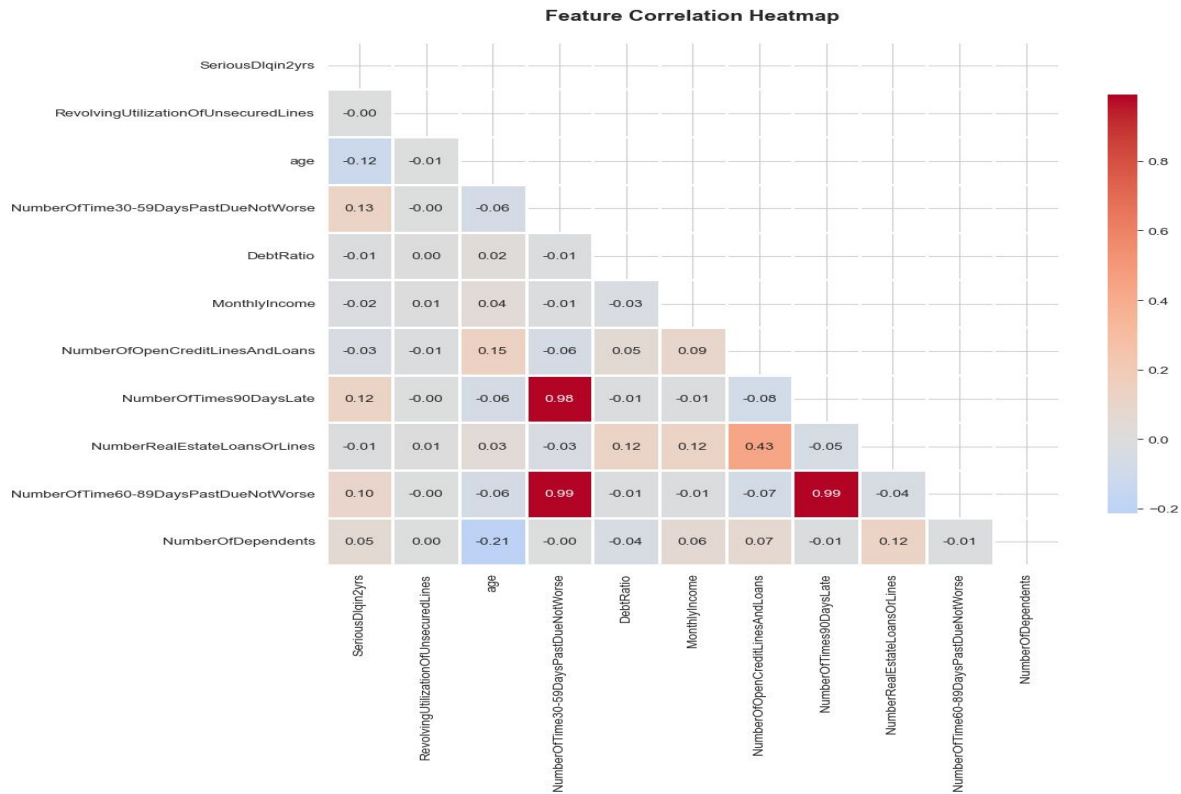
## Late Payment Behavior





# EDA

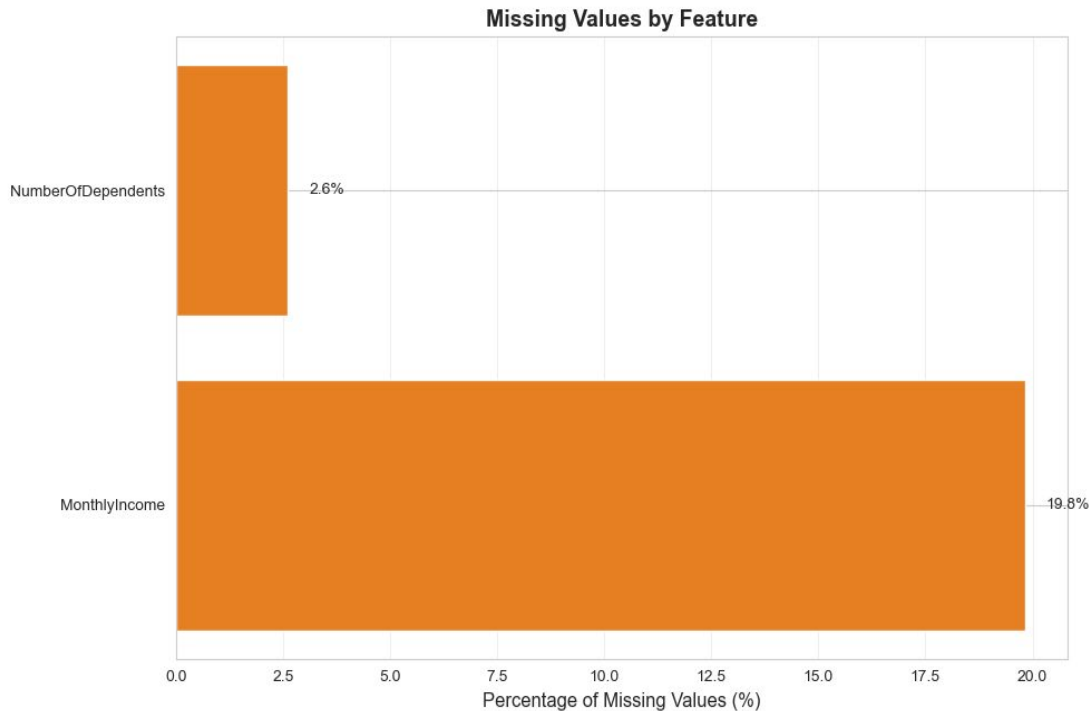
## Feature Correlation Analysis



# Data Preprocessing

## Missing Value Analysis

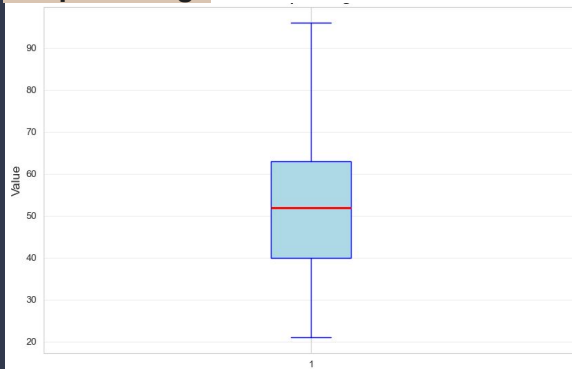
```
SeriousDlqin2yrs      0
RevolvingUtilizationOfUnsecuredLines  0
age                  0
NumberOfTime30-59DaysPastDueNotWorse  0
DebtRatio            0
MonthlyIncome        29731
NumberOfOpenCreditLinesAndLoans      0
NumberOfTimes90DaysLate              0
NumberRealEstateLoansOrLines         0
NumberOfTime60-89DaysPastDueNotWorse  0
NumberOfDependents    3924
dtype: int64
```



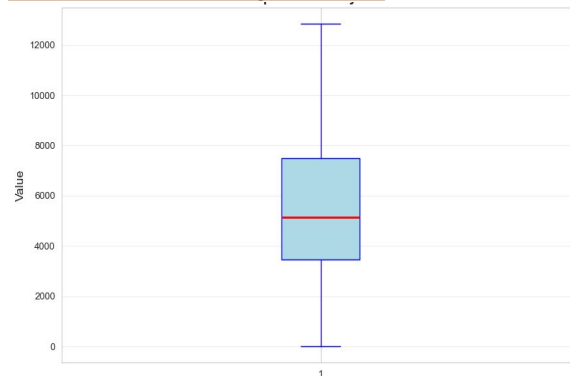
# Data Preprocessing

## Handling Outliers

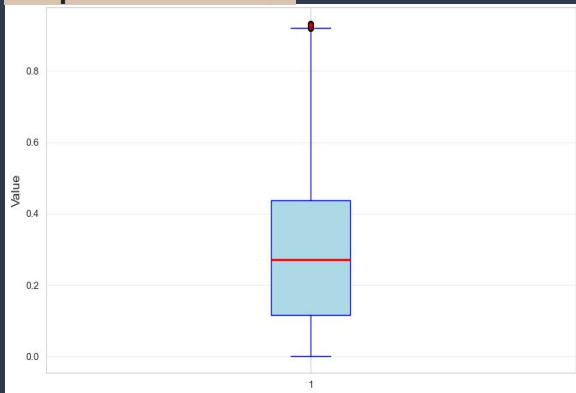
Boxplot of Age



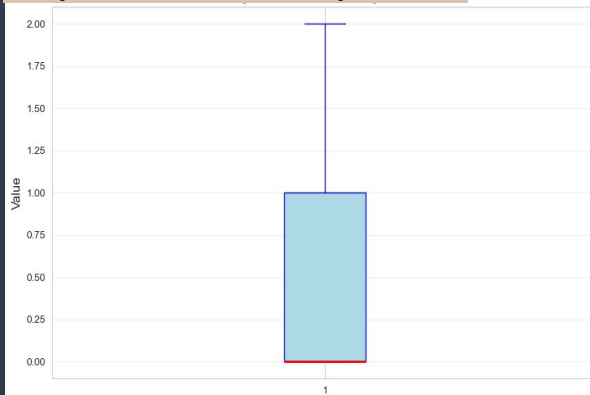
Boxplot of Monthly Income



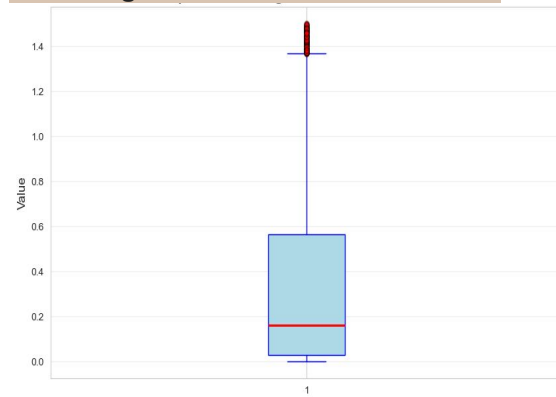
Boxplot of DebtRatio



Boxplot of NumberOfDependents



RevolvingUtilizationOfUnsecuredLines



# Model and Evaluation

- ✓ Logistic Regression
- ✓ Random Forest Classifier
- ✓ Gradient Boosting Classifier
- ✓ MLP Classifier

```
pipes = {}

for acronym, model in models.items():
    if acronym in ['lr', 'mlpc']:
        pipes[acronym] = Pipeline([
            ('scaler', StandardScaler()),
            ('smote', SMOTE(random_state=random_seed)),
            ('model', model),
        ])
    else:
        pipes[acronym] = Pipeline([
            ('smote', SMOTE(random_state=random_seed)),
            ('model', model),
        ])

pipes
```

# Hyperparameter Tuning

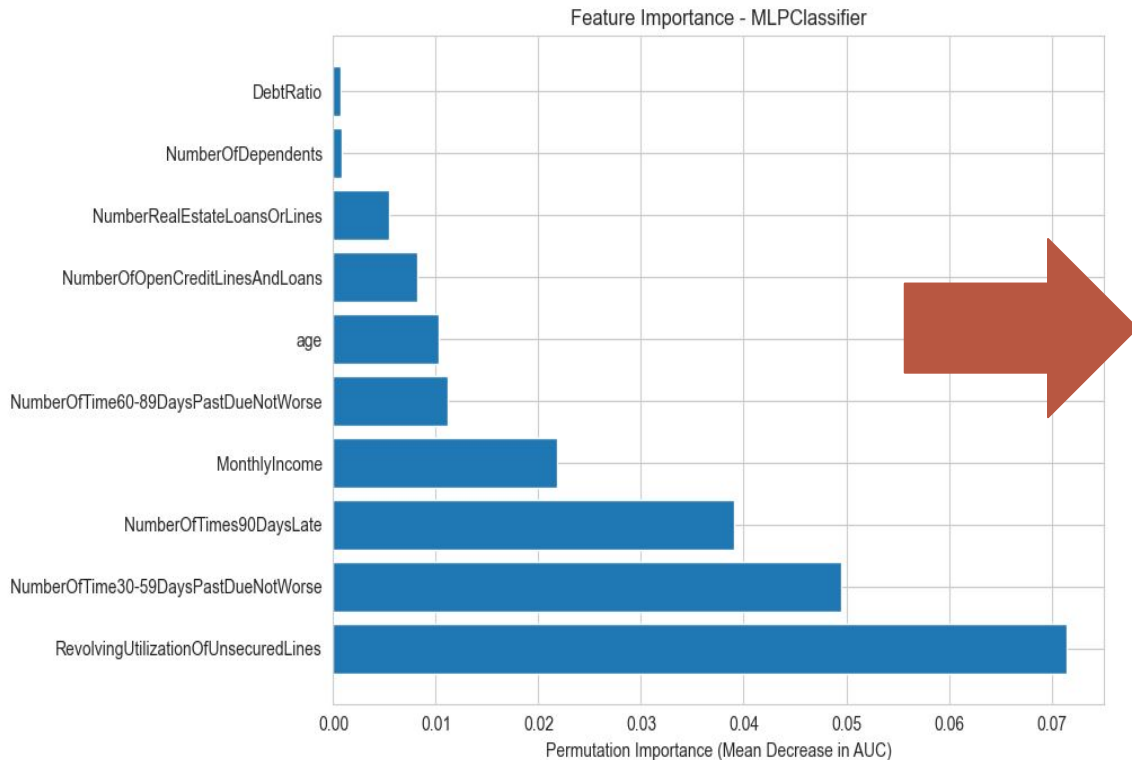
## GridSearch CV

Logistic Regression	<ul style="list-style-type: none"><li>• <code>model__C</code> : [0.1, 1, 10],</li><li>• <code>model__penalty</code> : ['l2'],</li><li>• <code>model__class_weight</code> : [None, 'balanced']</li></ul>
Random Forest Classifier	<ul style="list-style-type: none"><li>• <code>model__n_estimators</code> : [200, 400],</li><li>• <code>model__max_depth</code> : [5, 10, None],</li><li>• <code>model__min_samples_split</code> : [2, 5],</li><li>• <code>model__class_weight</code> : [None, 'balanced_subsample']</li></ul>
Gradient Boosting Classifier	<ul style="list-style-type: none"><li>• <code>model__n_estimators</code> : [100, 200],</li><li>• <code>model__learning_rate</code> : [0.05, 0.1],</li><li>• <code>model__max_depth</code> : [2, 3]</li></ul>
MLP Classifier	<ul style="list-style-type: none"><li>• <code>model__alpha</code> : <code>alpha_grids</code>,</li><li>• <code>model__learning_rate_init</code> : <code>lr_init_grids</code>,</li><li>• <code>model__hidden_layer_sizes</code> : <code>hidden_layer_sizes_grids</code></li></ul>

# Model Result

Model	best_cv_auc	test_auc
<b>MLP Classifier</b>	<b>0.851182</b>	<b>0.840149</b>
Random Forest Classifier	0.828728	0.821720
Gradient Boosting Classifier	0.826637	0.815682
Logistic Regression	0.800074	0.794796

# Feature Importance



## REDUCED FEATURE MODEL PERFORMANCE

Test AUC: 0.8190  
Test Accuracy: 0.9366  
Balanced Accuracy: 0.5753

Original Model AUC: 0.8495  
AUC Difference: -0.0305

## CLASSIFICATION REPORT:

	precision	recall	f1-score	support
No Delinquency	0.94	0.99	0.97	18559
Delinquency	0.62	0.16	0.25	1345
accuracy			0.94	19904
macro avg	0.78	0.58	0.61	19904
weighted avg	0.92	0.94	0.92	19904

## CONFUSION MATRIX:

```
[[18430  129]
 ...
 True Negatives: 18430
 False Positives: 129
 False Negatives: 1133
 True Positives: 212
```

Thank you

